



ISE

Industrial and
Systems Engineering

A service system with randomly behaving on-demand
agents

LAM M. NGUYEN AND ALEXANDER L. STOLYAR

Department of Industrial and Systems Engineering, Lehigh University, USA

ISE Technical Report 15T-016



A service system with randomly behaving on-demand agents

Lam M. Nguyen
Department of Industrial
and Systems Engineering
Lehigh University
Bethlehem, PA 18015
lmn214@lehigh.edu

Alexander L. Stolyar
Department of Industrial
and Systems Engineering
Lehigh University
Bethlehem, PA 18015
stolyar@lehigh.edu

November 30, 2015

Abstract

We consider a service system where agents (or, servers) are invited on-demand. Customers arrive as a Poisson process and join a customer queue. Customer service times are i.i.d. exponential. Agents' behavior is random in two respects. First, they can be invited into the system exogenously, and join the agent queue after a random time. Second, with some probability they rejoin the agent queue after a service completion, and otherwise leave the system. The customer and agent queues cannot be non-empty simultaneously – the head-of-the-line customer and agent are matched immediately and together go to service. The objective is to design a real-time adaptive agent invitation scheme that keeps both customer and agent queues/waiting-times small. We study an adaptive scheme, which controls the number of pending agent invitations, based on queue-state feedback.

We study the system process fluid limits, in the asymptotic regime where the customer arrival rate goes to infinity. The fluid limit trajectories have complicated behavior – there are two domains where they follow different ODE, and a “reflecting” boundary. We use the machinery of switched linear systems and Common Quadratic Lyapunov Functions to approach the stability of fluid limits at the desired equilibrium point (with zero queues). We derive sufficient local stability conditions for the fluid limits. When these conditions do hold, numerical and simulation experiments show good overall performance of the scheme. We conjecture that, for our model, local stability is in fact sufficient for global stability of fluid limits; the validity of this conjecture is supported by simulations.

1 Introduction

We study a service system with exogenously arriving customers, and servers, called *agents*, which can be invited to join the system at any time. The system control needs to match the arriving customers with invited agents, with the objective to minimize waiting times of both customers and agents. What makes this problem non-trivial is the fact that there is uncertainty in the agents'

behavior. First, invited agents do not arrive into the system immediately; instead they join the system after a random delay. Second, after an agent is done serving a customer, it can either leave the system or return to serve more customers.

This model (described in more detail below) is a generalization of that in [18, 13]. It was originally motivated by applications to call/contact centers (see e.g. [1, 12, 18] and references therein), where what we call agents are “special agents”, or “knowledge workers,” whose time is expensive, so that it is inefficient to have them working fixed shifts, with inevitable periods of idle time due to random fluctuations in customer demand. It is much more reasonable to invite them on-demand in real time; however, designing an efficient agent invitation strategy is non-trivial due to randomness in agent behavior. Besides efficiency (in terms of minimizing customer and agent waiting times), another highly desirable feature of the invitation scheme is simplicity and robustness.

We note that the model is generic and has other applications, or potential applications. One example is telemedicine [2], in which case “agents” are doctors, invited on-demand to serve patients remotely. Another example is crowdsourcing-based customer service [5, 3]. Also note that the model has relation to classical assemble-to-order models, where customers are orders and “invited agents” are products, which cannot be produced/assembled instantly. The model is also related to “double-ended queues” (see e.g. [9, 11]) and matching systems (see e.g. [6]); although in such models arrivals of all types into the system are typically exogenous, as opposed to being controlled.

More specifically, our model is as follows. Customers arrive as a Poisson process and join a customer queue. Customer service times are i.i.d. exponential. Agents’ behavior is random in two respects. First, they can be invited into the system exogenously, and join the agent queue after a random time. Second, with some probability they rejoin the agent queue after a service completion, and otherwise leave the system. (This generalizes the model in [18, 13], where the agents always leave the system after service completions, thus making our model more realistic in many scenarios.) The customer and agent queues cannot be non-empty simultaneously – the head-of-the-line customer and agent are matched immediately and together go to service. The objective is to design a real-time adaptive agent invitation scheme that keeps both customer and agent queues/waiting-times small.

We study a feedback-based adaptive scheme of [18, 13], which controls the number of pending agent invitations, depending on the customer and/or agent queue lengths and their changes. Due to the fact that our model is more general, the system dynamics is substantially more complicated.

The system state can be described by three variables, which are the number of pending invited agents, the difference between agent and customer queues, and the number of customers (or agents) in service. We consider an alternative, equivalent representation of the system state, which is also described by three variables: the number of pending invited agents, the difference between agent and customer queues, and the total number of customers and agents in the system.

Exact analysis of the model is prohibitively hard, and thus, we analyze it in an asymptotic regime where the customer arrival rate becomes large while the distribution of an agent response times and service times are fixed. We show convergence of the fluid-scaled process to the fluid limit (Theorem 1). The fluid limit trajectories have a complicated behavior – there are two domains where they follow different ODE, and a “reflecting” boundary. This presents challenges for proving stability of the fluid limits, namely the convergence of their trajectories to the equilibrium point, where the

queues are zero.

The focus of this paper and our **main results** concern the system *local stability* at the equilibrium point; specifically, the stability of the dynamic system which describes fluid limit trajectories when the state is away from the boundary. We use the machinery of switched linear systems and common quadratic Lyapunov functions [10, 17] to derive sufficient local stability conditions (Theorem 2). We conjecture that, for our model, local stability is in fact sufficient for global stability of fluid limits; the validity of this conjecture is supported by numerical and simulation experiments.

Our simulation experiments also show good overall performance of the feedback scheme when the local stability conditions hold.

1.1 Organization of the paper

The remainder of the paper is organized as follows. We will finish this section with basic notations, conventions, and abbreviations. Some necessary background facts on linear systems and switched linear systems are given in Section 2. In Section 3, we describe the model in detail. In Section 4 we state the main results of the paper. These results are proved in Sections 5 and 6. Simulation experiments are provided in Section 7; it also contains our conjectures about stability and local stability of fluid limits, supported by these simulations. We will conclude the paper in Section 8.

1.2 Basic notations, conventions and abbreviations

Sets of real and real non-negative numbers are denoted by \mathbb{R} and \mathbb{R}_+ ; \mathbb{R}^d and \mathbb{R}_+^d are the corresponding vector spaces. The standard Euclidean norm of a vector $x \in \mathbb{R}^n$ is denoted $\|x\|$. For a vector a or matrix A , we write their transposes as a^T or A^T . We write $x(\cdot)$ to mean the function (or random process) $(x(t), t \geq 0)$. For a real-valued function $x(\cdot) : \mathbb{R}_+ \rightarrow \mathbb{R}$, we use either $x'(t)$ or $(d/dt)x(t)$ to denote the derivative with respect to t , and for $x(\cdot) : \mathbb{R}_+ \rightarrow \mathbb{R}^d$, we write $(d/dt)x(t) = (x'_1(t), \dots, x'_d(t))$. For a real number x , let $x^+ = \max\{x, 0\}$ and $x^- = -\min\{x, 0\}$ and let

$$\text{sgn}(x) = \begin{cases} 1, & x > 0 \\ 0, & x = 0 \\ -1, & x < 0 \end{cases}$$

For $x, y \in \mathbb{R}$, we denote $x \wedge y = \min\{x, y\}$ and $x \vee y = \max\{x, y\}$. Symbol \Leftrightarrow means “equivalent to”. We write $x^r \rightarrow x \in \mathbb{R}^n$ to denote ordinary convergence in \mathbb{R}^n . For a finite set of scalar functions $f_n(t), t \geq 0, n \in \mathbb{N}$, a point t is called *regular* if for any subset $\mathbb{N}_0 \subseteq \mathbb{N}$, the derivatives

$$\frac{d}{dt} \max_{n \in \mathbb{N}_0} f_n(t) \text{ and } \frac{d}{dt} \min_{n \in \mathbb{N}_0} f_n(t)$$

exist. (To be precise, we require that each derivative is proper: both left and right derivatives exist and are equal.) We use small- o notation for deterministic function: for two real-valued functions f and g , we write $f(x) = o(g(x))$ if $\limsup_{x \rightarrow \infty} |f(x)/g(x)| = 0$.

Abbreviation *u.o.c.* means *uniform on compact sets* convergence of functions, with the argument determined by the context (usually in $[0, \infty)$); *w.p.1* means *with probability 1*; *i.i.d.* means *independent identically distributed*; RHS means *right hand side*; FLLN means *functional strong law of large numbers*; CQLF means *common quadratic Lyapunov function*.

2 Necessary background facts

2.1 Definitions and results related to switched linear system

In this paper, we will use some machinery of switched linear system. Here, we provide some necessary background. Consider a *switched linear system*

$$\Sigma_S : u'(t) = A(t)u(t) , A(t) \in \mathcal{A} = \{A_1, \dots, A_m\} \quad (1)$$

where \mathcal{A} is a set of matrices in $\mathbb{R}^{n \times n}$, and $t \rightarrow A(t)$ is a piecewise constant mapping from nonnegative real numbers into \mathcal{A} . (Note that by definition in [17] piecewise constant maps have only finitely many discontinuities in any bounded time-interval.) For $1 \leq i \leq m$, the i^{th} constituent system of the switched linear system (1) is the *linear time-invariant (LTI) system*

$$\Sigma_{A_i} : u'(t) = A_i u(t). \quad (2)$$

The origin is an *exponentially stable equilibrium* of a switched linear system Σ_s if there exists real constants $C > 0$, $a > 0$ such that $\|u(t)\| \leq Ce^{-at}\|u(0)\|$ for $t \geq 0$, for all solutions $u(t)$ of the system (1) under any $A(t)$. (see [7, 17])

A symmetric square $n \times n$ matrix M with real coefficients is *positive definite* if $z^T M z > 0$ for every non-zero column vector $z \in \mathbb{R}^n$. A symmetric square $n \times n$ matrix M with real coefficients is *negative definite* if $z^T M z < 0$ for every non-zero column vector $z \in \mathbb{R}^n$. A square matrix A is called a *Hurwitz matrix* (or *stable matrix*) if every eigenvalue of A has strictly negative real part. (see [15])

The function $V(u) = u^T P u$ is a *quadratic Lyapunov function* (QLF) for the system $\Sigma_A : u'(t) = A u(t)$ if (i) P is symmetric and positive definite, and (ii) $PA + A^T P$ is negative definite. Let $\{A_1, \dots, A_m\}$ be a collection of $n \times n$ Hurwitz matrices, with associated stable LTI systems $\Sigma_{A_1}, \dots, \Sigma_{A_m}$. Then the function $V(u) = u^T P u$ is a *common quadratic Lyapunov function* (CQLF) for these systems if V is a QLF for each individual system. (see [10, 17])

The following propositions will be used in the proof of our result (Theorem 2).

Proposition 1. (see [10, 17]). *The existence of a CQLF for the LTI systems is sufficient for the exponential stability of a switched linear system.*

Proposition 2. (see [10, 17]). *Let A^+ and A^- be Hurwitz matrices in $\mathbb{R}^{n \times n}$, and the difference $A^+ - A^-$ has rank one. Then two systems $u'(t) = A^+ u(t)$ and $u'(t) = A^- u(t)$ have a CQLF if and only if the matrix product $A^+ A^-$ has no negative real eigenvalues.*

2.2 Stability of linear systems

The following propositions will also be used in the proof of our result (Theorem 2).

Proposition 3. (Routh-Hurwitz stability criterion) (see [15]). *Let $\det(A - \lambda I) = 0$ be the characteristic equation of matrix A , which is equivalent to*

$$L(\lambda) = a_0 \lambda^3 + a_1 \lambda^2 + a_2 \lambda + a_3 = 0 , a_0 > 0. \quad (3)$$

Matrix A is Hurwitz if and only if a_1, a_2, a_3 are positive and satisfying $a_1 a_2 > a_0 a_3$.

Proposition 4. (see [8]). *The general cubic equation has the form*

$$a\lambda^3 + b\lambda^2 + c\lambda + d = 0, \quad a \neq 0, \quad (4)$$

and discriminant

$$\Delta = 18abcd - 4b^3d + b^2c^2 - 4ac^3 - 27a^2d^2. \quad (5)$$

If $\Delta > 0$, then the equation has three distinct real roots.

If $\Delta = 0$, then the equation has a multiple root and all its roots are real.

If $\Delta < 0$, then the equation has one real root and two nonreal complex conjugate roots.

3 Model and algorithm

Our model is a generalization of that considered in [18, 13]. Customers arrive according to a Poisson process of rate $\Lambda > 0$, and join the customer queue waiting for an available agent and are served in the order of their arrival. There is an infinite pool of potential agents, which can be invited to serve customers. Once being invited, an agent will respond after an independent exponentially distributed random time, with mean $1/\tilde{\beta}$; it accepts the invitation with probability $a > 0$, and otherwise rejects it. Let $\beta = a\tilde{\beta} > 0$ be the rate at which an agent accepts the invitation. Agents who accept their invitations join the agent queue, in the order of their arrival. The customer and agent queues cannot be positive simultaneously: the head-of-the-line customer and agents are immediately matched, leave their queues, and together go to service. Each service time is an exponentially distributed random variable with mean $1/\mu$; after the service completion, the customer leaves the system, while the agent rejoins the agent queue with probability $\alpha \in [0, 1)$. Thus, there are two ways in which agents join the queue – exogenously invited agents accepting invitations and agents already in the system rejoining the queue after service completions. (The model in [18, 13] is a special case of ours, with $\alpha = 0$; in other words, the agents certainly leave the system after service completions, and therefore there is no need to account for agents being in service.)

Let $X(t)$ be the number of pending agents that have been invited but have not decided to accept or decline the invitations at time t . Let $Q_c(t)$ be the number of customers in the customer queue at time t . Let $Q_a(t)$ be the number of agents in the agent queue at time t . And we also define $Y(t) = Q_a(t) - Q_c(t)$ as the difference of the agent queue and customer queue at time t . Let $Z(t)$ be the number of customers (or agents) in service at time t . We assume that the non-idling condition holds, that is, agents do not idle when there are customers waiting in the customer queue, which means that at each time t , either the customer queue or the agent queue must be empty. The system state can be described by three variables: X : 'the number of pending invited agents'. Y : 'the difference between agent and customer queues'. Z : 'the number of customers (or agents) in service'. Figure 1 depicts such an agent invitation system.

The feedback invitation scheme in [18], let us label it as *Scheme A*, is defined as follows. The scheme maintains a ‘‘target’’ $X_{target}(t)$ for the number of invited agents $X(t)$. The target $X_{target}(t)$ is changed by $\Delta X_{target}(t) = [-\gamma\Delta Y(t) - \epsilon Y(t)\Delta t]$ at each time t when $Y(t)$ changes by $\Delta Y(t)$ (which can be either $+1$ or -1), where $\gamma > 0$ and $\epsilon > 0$ are the algorithm parameters and Δt is the time duration from the previous change of Y . New agents are invited if and only if $X(t) < X_{target}(t)$, where $X(t)$ is the actual number of invited (pending) agents; therefore, $X(t) \geq X_{target}(t)$ holds at all times. In addition, the target $X_{target}(t)$ is not allowed to go below zero, $X_{target}(t) \geq 0$; i.e. if an

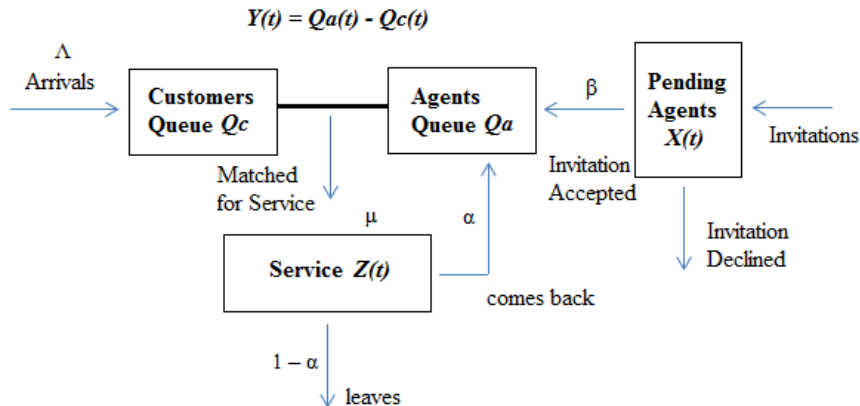


Figure 1: An Agent Invitation System

update of $X_{target}(t)$ makes it negative, its value is immediately reset to zero. Note that $X_{target}(t)$ is not necessarily an integer.

Although the scheme we consider is same as in [18], the model we apply it to is different. Namely, arrivals into the agent queue are not only due to invited agents accepting invitations, but also due to agents returning immediately after the service completions. As a result, the process describing the system evolution contains additional variable Z , and is more complicated.

To simplify our theoretical analysis, just as in [13], we consider a “stylized” version of Scheme A, which has the same basic dynamics, but keeps $X_{target}(t)$ integer and assumes that $X(t) = X_{target}(t)$ at all times; the latter is equivalent to assuming that not only agent invitations can be issued instantly, but they can also be withdrawn at any time. Given these assumptions, when pending agents decline invitations, it has no impact on the system state, because $X(t)$ is immediately “replenished” by inviting another agent. Therefore, in the analysis of stylized scheme, the events of declined invitations can be ignored.

Formally, the stylized scheme, which we label *Scheme B*, is defined as follows. There are four types of mutually independent, and independent of the past, events that affect the dynamics of $X(t)$, $Y(t)$ and $Z(t)$ in a small time interval $[t, t + dt]$: (i) a customer arrival with probability $\Lambda dt + o(dt)$, (ii) an agent acceptance with probability $\beta X(t) dt + o(dt)$, (iii) an additional event with probability $\epsilon |Y(t)| dt + o(dt)$, and (iv) service completion with probability $\mu Z(t) dt + o(dt)$.

The changes at these event times are described as follows:

(i) Upon a customer arrival, if $Y(t) > 0$, $Z(t)$ changes by $\Delta Z(t) = 1$; and if $Y(t) \leq 0$, $Z(t)$ changes by $\Delta Z(t) = 0$. $Y(t)$ changes by $\Delta Y(t) = -1$, and $X(t)$ changes by $\Delta X(t) = \gamma$ (we assume that $\gamma > 0$ is an integer).

(ii) Upon the acceptance of an invitation, if $Y(t) < 0$, $Z(t)$ changes by $\Delta Z(t) = 1$; and if $Y(t) \geq 0$, $Z(t)$ changes by $\Delta Z(t) = 0$. $Y(t)$ changes by $\Delta Y(t) = 1$, and $X(t)$ changes by $\Delta X(t) = -(\gamma \wedge X(t))$, that is, the change is by $-\gamma$ but $X(t)$ is kept to be nonnegative.

(iii) Upon the third type of event, if $X(t) \geq 1$, the change $\Delta X(t) = -\text{sgn}(Y(t))$ occurs; and if $X(t) = 0$, the change $\Delta X(t) = 1$ occurs if $Y(t) < 0$ and $\Delta X(t) = 0$ if $Y(t) \geq 0$.

(iv) Upon the service completion, (a) with probability α , if $Y(t) < 0$, the change $\Delta Z(t) = -1 + 1 = 0$ occurs; and if $Y(t) \geq 0$, the change $\Delta Z(t) = -1$ occurs; $Y(t)$ changes by $\Delta Y(t) = 1$, and

$\Delta X(t) = -(\gamma \wedge X(t))$. (b) With probability $(1 - \alpha)$, $Z(t)$ changes by $\Delta Z(t) = -1$.

4 Main Results

We consider a sequence of systems, indexed by a scaling parameter $r \rightarrow \infty$. In the system with index r , the arrival rate is λr , while the parameters $\alpha, \beta, \mu, \epsilon, \gamma$ are constant. The corresponding process is (X^r, Y^r, Z^r) , where $X^r = (X^r(t), t \geq 0)$, $Y^r = (Y^r(t), t \geq 0)$ and $Z^r = (Z^r(t), t \geq 0)$. We will center the values of X^r, Y^r , and Z^r by $\lambda r(1 - \alpha)/\beta, 0$, and $\lambda r/\mu$, respectively. These values are such that $\beta X^r + \mu \alpha Z^r = \lambda r$, which means that on average the arrival rate of agents into the agent queue matches the rate of customer arrivals. We define fluid-scaled processes with centering

$$\begin{cases} \bar{X}^r = \frac{1}{r} \left(X^r - \frac{\lambda r(1-\alpha)}{\beta} \right) \\ \bar{Y}^r = \frac{1}{r} Y^r \\ \bar{Z}^r = \frac{1}{r} \left(Z^r - \frac{\lambda r}{\mu} \right). \end{cases} \quad (6)$$

Let W be the total number of customers and agents in the system. We know that Y is the difference between agent and customer queues (only one of those queues can be positive at any time since we have the non-idling condition) and Z is the number of customers (or agents) in service. From this, $W = |Y| + 2Z$, which is equivalent to $Z = \frac{1}{2}(W - |Y|)$. Instead of using the process (X, Y, Z) , we are using a new process (X, Y, W) . This process (X, Y, W) is more convenient for the analysis. We have new fluid-scaled processes with centering

$$\begin{cases} \bar{X}^r = \frac{1}{r} \left(X^r - \frac{\lambda r(1-\alpha)}{\beta} \right) \\ \bar{Y}^r = \frac{1}{r} Y^r \\ \bar{W}^r = \frac{1}{r} \left(W^r - \frac{2\lambda r}{\mu} \right). \end{cases} \quad (7)$$

Theorem 1. *Consider a sequence of processes $(\bar{X}^r, \bar{Y}^r, \bar{W}^r)$, $r \rightarrow \infty$, with deterministic initial states such that $(\bar{X}^r(0), \bar{Y}^r(0), \bar{W}^r(0)) \rightarrow (x(0), y(0), w(0))$ for some fixed $(x(0), y(0), w(0)) \in \mathbb{R}^3$, $x(0) \geq -\frac{\lambda(1-\alpha)}{\beta}$. Then, these processes can be constructed on a common probability space, so that the following holds. W.p.1, from any subsequence of r , there exists a further subsequence such that*

$$(\bar{X}^r, \bar{Y}^r, \bar{W}^r) \rightarrow (x, y, w) \text{ u.o.c. as } r \rightarrow \infty \quad (8)$$

where (x, y, w) is a locally Lipschitz trajectory such that at any regular point $t \geq 0$

$$\begin{cases} x'(t) = \begin{cases} -\gamma y'(t) - \epsilon y, & \text{if } x(t) > -\frac{\lambda(1-\alpha)}{\beta} \\ [-\gamma y'(t) - \epsilon y] \vee 0, & \text{if } x(t) = -\frac{\lambda(1-\alpha)}{\beta} \end{cases} \\ y'(t) = \beta x + \frac{1}{2} \alpha \mu (w - |y|) \\ w'(t) = \beta x + \frac{1}{2} (\alpha - 2) \mu (w - |y|). \end{cases} \quad (9)$$

A limit trajectory (x, y, w) specified in Theorem 1 will be called a *fluid limit* starting from $(x(0), y(0), w(0))$.

Consider a dynamic system $(x(t), y(t), w(t)) \in \mathbb{R}^3$:

$$\begin{cases} x'(t) = -\gamma y'(t) - \epsilon y \\ y'(t) = \beta x + \frac{1}{2} \alpha \mu (w - |y|) \\ w'(t) = \beta x + \frac{1}{2} (\alpha - 2) \mu (w - |y|). \end{cases} \quad (10)$$

This dynamic system describes the dynamics of fluid limit trajectories when the state is away from the boundary $x = -\frac{\lambda(1-\alpha)}{\beta}$. The *non-linear* system (10) is a generalization of the linear system, considered in [13]. The latter is a special case of (10) without variable w , and with $\alpha = 0$. The system in [13] is simply linear, while (10) has two domains, defined by the sign of y . The following result provides sufficient exponential stability conditions for the system (10).

Theorem 2. (Sufficient exponential stability conditions). *For any set of positive β , μ , and $\alpha \in (0, 1)$, there exist values of $\gamma > 0$ and $\epsilon > 0$ satisfying the following conditions*

$$\begin{cases} \frac{\beta\gamma^2}{4} < \epsilon < \frac{\beta\gamma^2}{2} \\ \epsilon > \frac{\beta\gamma^2}{2} - \left(\frac{\alpha\gamma\mu}{2} - \frac{(1-\alpha)\mu^2}{2\beta} \right) \\ \gamma > \frac{(1-\alpha)\mu}{\alpha\beta}. \end{cases} \quad (11)$$

For the parameters, satisfying these conditions, common quadratic Lyapunov function (CQLF) of the system (10) exists, and the system (10) is exponentially stable.

We say that our fluid-limit system is *globally stable* if every fluid limit trajectory converges to the equilibrium point $(0, 0, 0)$; we say that it is *locally stable* if every trajectory of the dynamic system (10) converges to the equilibrium point $(0, 0, 0)$. Therefore, the conditions (11) in Theorem 2 are also sufficient for the local stability of our system.

5 Fluid scale analysis and proof of Theorem 1

In order to prove Theorem 1, it suffices to show that w.p.1 from any subsequence of r , we can choose a further subsequence, along which a u.o.c. convergence to a fluid limit holds.

Given the initial state $(X^r(0), Y^r(0), W^r(0))$, we construct the process (X^r, Y^r, W^r) , for all r , on the same probability space via a common set of independent Poisson process as follows:

$$X^r(t) = G^r(t) + \left(- \min_{0 \leq s \leq t} G^r(s) \right) \vee 0, \quad (12)$$

$$\begin{aligned} G^r(t) = & X^r(0) + \gamma N_1(\lambda r t) - \gamma N_2 \left(\beta \int_0^t X^r(s) ds \right) - \gamma N_4 \left(\alpha \mu \int_0^t \frac{1}{2} (W^r(s) - |Y^r(s)|) ds \right) + \\ & + N_5 \left(\epsilon \int_0^t (Y^r(s))^- ds \right) - N_6 \left(\epsilon \int_0^t (Y^r(s))^+ ds \right), \end{aligned} \quad (13)$$

$$Y^r(t) = Y^r(0) + N_2 \left(\beta \int_0^t X^r(s) ds \right) + N_4 \left(\alpha \mu \int_0^t \frac{1}{2} (W^r(s) - |Y^r(s)|) ds \right) - N_1(\lambda r t), \quad (14)$$

$$\begin{aligned} W^r(t) = & W^r(0) + N_1(\lambda r t) + N_2 \left(\int_0^t \beta X^r(s) ds \right) - N_3 \left(\int_0^t 2(1-\alpha)\mu \frac{1}{2} (W^r(s) - |Y^r(s)|) ds \right) - \\ & - N_4 \left(\int_0^t \alpha \mu \frac{1}{2} (W^r(s) - |Y^r(s)|) ds \right), \end{aligned} \quad (15)$$

and $N_i(\cdot)$, $i = 1, \dots, 6$ are mutually independent unit-rate Poisson processes [14]. N_1 is the process which drives customer arrivals. N_2 is the process which drives the acceptance of invitations. N_3 is the process which drives the service completions, with agents leaving the system. N_4 is the process which drives the service completions, with agents coming back. N_5 and N_6 are the processes which drive the third type of event. W.p.1, for any r , relations (12)-(15) uniquely define the realization

of (X^r, Y^r, W^r) via the realizations of the driving processes $N_i(\cdot)$. Relation (12), the “reflection” at zero, corresponds to the property that $X^r(t)$ cannot become negative.

The functional strong law of large numbers (FSLLN) holds for each Poisson process N_i :

$$\frac{N_i(rt)}{r} \rightarrow t, \quad r \rightarrow \infty, \quad \text{u.o.c., w.p.1.} \quad (16)$$

We consider the sequence of associated fluid-scaled processes $(\bar{X}^r, \bar{Y}^r, \bar{W}^r)$ as defined in (7). (Note that these processes are centered.) Let a constant $m > \|(x(0), y(0), w(0))\|$ be fixed. For each r , on the same probability space as $(\bar{X}^r, \bar{Y}^r, \bar{W}^r)$, let us define a modified fluid-scaled process $(\bar{X}_m^r, \bar{Y}_m^r, \bar{W}_m^r)$. Let $(\bar{X}_m^r, \bar{Y}_m^r, \bar{W}_m^r)$ start from the same initial state as $(\bar{X}^r, \bar{Y}^r, \bar{W}^r)$, i.e., $(\bar{X}_m^r(0), \bar{Y}_m^r(0), \bar{W}_m^r(0)) = (\bar{X}^r(0), \bar{Y}^r(0), \bar{W}^r(0))$. The modified process $(\bar{X}_m^r, \bar{Y}_m^r, \bar{W}_m^r)$ follows the same path as $(\bar{X}^r, \bar{Y}^r, \bar{W}^r)$ until the first time that $\|(\bar{X}^r(t), \bar{Y}^r(t), \bar{W}^r(t))\| \geq m$. Denote this time by τ_m^r . We then freeze the process $(\bar{X}_m^r, \bar{Y}_m^r, \bar{W}_m^r)$ at the value $(\bar{X}^r(\tau_m^r), \bar{Y}^r(\tau_m^r), \bar{W}^r(\tau_m^r))$, i.e. $(\bar{X}_m^r(t), \bar{Y}_m^r(t), \bar{W}_m^r(t)) = (\bar{X}^r(\tau_m^r), \bar{Y}^r(\tau_m^r), \bar{W}^r(\tau_m^r))$ for all $t \geq \tau_m^r$.

Lemma 1. *Fix $(x(0), y(0), w(0))$ and a finite constant $m > \|(x(0), y(0), w(0))\|$. Then, w.p.1 for any subsequence of r , there exists a further subsequence, along which $(\bar{X}_m^r, \bar{Y}_m^r, \bar{W}_m^r)$ converges u.o.c. to a Lipschitz continuous trajectory (x_m, y_m, w_m) , which satisfies properties (9) at any regular time $t \geq 0$ such that $\|(x_m(t), y_m(t), w_m(t))\| < m$.*

Proof. For the modified fluid-scaled processes $(\bar{X}_m^r, \bar{Y}_m^r, \bar{W}_m^r)$, we define the associated counting processes for upward and downward jumps. For $t \leq \tau_m^r$,

$$\bar{X}_m^{r\uparrow}(t) = r^{-1}\gamma N_1(\lambda rt) + r^{-1}N_5\left(\epsilon r \int_0^t (\bar{Y}_m^r(s))^- ds\right), \quad (17)$$

$$\begin{aligned} \bar{X}_m^{r\downarrow}(t) &= r^{-1}\gamma N_2\left(\beta r \int_0^t \left[\bar{X}_m^r(s) + \frac{\lambda(1-\alpha)}{\beta}\right] ds\right) + \\ &+ r^{-1}\gamma N_4\left(\frac{1}{2}\alpha\mu r \int_0^t \left[\bar{W}_m^r(s) + \frac{2\lambda}{\mu} - |\bar{Y}_m^r(s)|\right] ds\right) + r^{-1}N_6\left(\epsilon r \int_0^t (\bar{Y}_m^r(s))^+ ds\right), \end{aligned} \quad (18)$$

$$\bar{Y}_m^{r\uparrow}(t) = r^{-1}N_2\left(\beta r \int_0^t \left[\bar{X}_m^r(s) + \frac{\lambda(1-\alpha)}{\beta}\right] ds\right) + r^{-1}N_4\left(\frac{1}{2}\alpha\mu r \int_0^t \left[\bar{W}_m^r(s) + \frac{2\lambda}{\mu} - |\bar{Y}_m^r(s)|\right] ds\right), \quad (19)$$

$$\bar{Y}_m^{r\downarrow}(t) = r^{-1}N_1(\lambda rt), \quad (20)$$

$$\bar{W}_m^{r\uparrow}(t) = r^{-1}N_1(\lambda rt) + r^{-1}N_2\left(\beta r \int_0^t \left[\bar{X}_m^r(s) + \frac{\lambda(1-\alpha)}{\beta}\right] ds\right), \quad (21)$$

$$\begin{aligned} \bar{W}_m^{r\downarrow}(t) &= r^{-1}N_3\left((1-\alpha)\mu r \int_0^t \left[\bar{W}_m^r(s) + \frac{2\lambda}{\mu} - |\bar{Y}_m^r(s)|\right] ds\right) + \\ &+ r^{-1}N_4\left(\frac{1}{2}\alpha\mu r \int_0^t \left[\bar{W}_m^r(s) + \frac{2\lambda}{\mu} - |\bar{Y}_m^r(s)|\right] ds\right), \end{aligned} \quad (22)$$

and for $t > \tau_m^r$, all these counting processes are frozen at their values at time τ_m^r , that is,

$$\begin{cases} \bar{X}_m^{r\uparrow}(t) = \bar{X}_m^{r\uparrow}(\tau_m^r), & \bar{X}_m^{r\downarrow}(t) = \bar{X}_m^{r\downarrow}(\tau_m^r), \\ \bar{Y}_m^{r\uparrow}(t) = \bar{Y}_m^{r\uparrow}(\tau_m^r), & \bar{Y}_m^{r\downarrow}(t) = \bar{Y}_m^{r\downarrow}(\tau_m^r), \\ \bar{W}_m^{r\uparrow}(t) = \bar{W}_m^{r\uparrow}(\tau_m^r), & \bar{W}_m^{r\downarrow}(t) = \bar{W}_m^{r\downarrow}(\tau_m^r). \end{cases} \quad (23)$$

Using the relations (12)-(15) and the fact that for $0 \leq t \leq \tau_m^r$ the original process $(\bar{X}^r, \bar{Y}^r, \bar{W}^r)$ and the modified process $(\bar{X}_m^r, \bar{Y}_m^r, \bar{W}_m^r)$ coincide, we have for all $t \geq 0$,

$$\bar{X}_m^r(t) = \bar{G}_m^r(t) + \left(-\lambda(1-\alpha)/\beta - \min_{0 \leq s \leq t} \bar{G}_m^r(s) \right) \vee 0, \quad (24)$$

$$\bar{G}_m^r(t) = \bar{X}^r(0) + \bar{X}_m^{r\uparrow}(t) - \bar{X}_m^{r\downarrow}(t), \quad (25)$$

$$\bar{Y}_m^r(t) = \bar{Y}^r(0) + \bar{Y}_m^{r\uparrow}(t) - \bar{Y}_m^{r\downarrow}(t), \quad (26)$$

$$\bar{W}_m^r(t) = \bar{W}^r(0) + \bar{W}_m^{r\uparrow}(t) - \bar{W}_m^{r\downarrow}(t). \quad (27)$$

The counting processes $\bar{X}_m^{r\uparrow}, \bar{X}_m^{r\downarrow}, \bar{Y}_m^{r\uparrow}, \bar{Y}_m^{r\downarrow}, \bar{W}_m^{r\uparrow}, \bar{W}_m^{r\downarrow}$ are non-decreasing. Using the Functional Strong Law of Large Number (FSLLN) (16) and the fact that the processes \bar{X}_m^r, \bar{Y}_m^r , and \bar{W}_m^r are uniformly bounded by construction, we see that w.p.1. for any subsequence of r , there exists a further subsequence along which the set of trajectories $(\bar{X}_m^{r\uparrow}, \bar{X}_m^{r\downarrow}, \bar{Y}_m^{r\uparrow}, \bar{Y}_m^{r\downarrow}, \bar{W}_m^{r\uparrow}, \bar{W}_m^{r\downarrow})$ converges u.o.c. to a set of non-decreasing Lipschitz continuous functions $(x_m^\uparrow, x_m^\downarrow, y_m^\uparrow, y_m^\downarrow, w_m^\uparrow, w_m^\downarrow)$. But then the u.o.c. convergence of $(\bar{X}_m^r, \bar{Y}_m^r, \bar{W}_m^r, \bar{G}_m^r)$ to a set of Lipschitz continuous functions (x_m, y_m, w_m, g_m) holds, where

$$x_m(t) = g_m(t) + \left(-\lambda(1-\alpha)/\beta - \min_{0 \leq s \leq t} g_m(s) \right) \vee 0, \quad (28)$$

$$g_m(t) = x(0) + x_m^\uparrow(t) - x_m^\downarrow(t), \quad (29)$$

$$y_m(t) = y(0) + y_m^\uparrow(t) - y_m^\downarrow(t), \quad (30)$$

$$w_m(t) = w(0) + w_m^\uparrow(t) - w_m^\downarrow(t), \quad (31)$$

and the following holds for t before fluid trajectory hits $\|(x_m(t), y_m(t), w_m(t))\| = m$

$$x_m^\uparrow(t) = \gamma\lambda t + \epsilon \int_0^t y_m^-(s) ds, \quad (32)$$

$$x_m^\downarrow(t) = \gamma\beta \int_0^t \left(x_m(s) + \frac{\lambda(1-\alpha)}{\beta} \right) ds + \frac{1}{2}\gamma\alpha\mu \int_0^t \left(w_m(s) + \frac{2\lambda}{\mu} - |y_m(s)| \right) ds + \epsilon \int_0^t y_m^+(s) ds, \quad (33)$$

$$y_m^\uparrow(t) = \beta \int_0^t \left(x_m(s) + \frac{\lambda(1-\alpha)}{\beta} \right) ds + \frac{1}{2}\alpha\mu \int_0^t \left(w_m(s) + \frac{2\lambda}{\mu} - |y_m(s)| \right) ds, \quad (34)$$

$$y_m^\downarrow(t) = \lambda t, \quad (35)$$

$$w_m^\uparrow(t) = \lambda t + \beta \int_0^t \left(x_m(s) + \frac{\lambda(1-\alpha)}{\beta} \right) ds, \quad (36)$$

$$w_m^\downarrow(t) = (1-\alpha)\mu \int_0^t \left(w_m(s) + \frac{2\lambda}{\mu} - |y_m(s)| \right) ds + \frac{1}{2}\alpha\mu \int_0^t \left(w_m(s) + \frac{2\lambda}{\mu} - |y_m(s)| \right) ds. \quad (37)$$

It is easy to verify that, for t before fluid trajectory hits $\|(x_m(t), y_m(t), w_m(t))\| = m$

$$\begin{cases} x'_m(t) = \begin{cases} -\gamma\beta x_m - \frac{1}{2}\gamma\alpha\mu w_m + \frac{1}{2}\gamma\alpha\mu |y_m| - \epsilon y_m, & \text{if } x_m(t) > -\frac{\lambda(1-\alpha)}{\beta} \\ [-\gamma\beta x_m - \frac{1}{2}\gamma\alpha\mu w_m + \frac{1}{2}\gamma\alpha\mu |y_m| - \epsilon y_m] \vee 0, & \text{if } x_m(t) = -\frac{\lambda(1-\alpha)}{\beta} \end{cases} \\ y'_m(t) = \beta x_m + \frac{1}{2}\alpha\mu(w_m - |y_m|) \\ w'_m(t) = \beta x_m + \frac{1}{2}(\alpha - 2)\mu(w_m - |y_m|) \end{cases} \quad (38)$$

which is equivalent to

$$\begin{cases} x'_m(t) = \begin{cases} -\gamma y'_m(t) - \epsilon y_m, & \text{if } x_m(t) > -\frac{\lambda(1-\alpha)}{\beta} \\ [-\gamma y'_m(t) - \epsilon y_m] \vee 0, & \text{if } x_m(t) = -\frac{\lambda(1-\alpha)}{\beta} \end{cases} \\ y'_m(t) = \beta x_m + \frac{1}{2}\alpha\mu(w_m - |y_m|) \\ w'_m(t) = \beta x_m + \frac{1}{2}(\alpha - 2)\mu(w_m - |y_m|). \end{cases} \quad (39)$$

This means properties (9) hold for the trajectory (x_m, y_m, w_m) . This completes the proof. \square

Conclusion of the proof of Theorem 1. It is obvious that inequality $\frac{d}{dt}\|(x_m(t), y_m(t), w_m(t))\| \leq C\|(x_m(t), y_m(t), w_m(t))\|$ holds for any m , and some common $C > 0$. From Gronwall's inequality [4], we have $\|(x_m(t), y_m(t), w_m(t))\| \leq \|(x(0), y(0), w(0))\|e^{Ct}$ for $t \geq 0$. For a given $(x(0), y(0), w(0))$, let us fix $T_l > 0$ and choose $m_l > \|(x(0), y(0), w(0))\|e^{CT_l}$. For this $T_l > 0$, there exists a subsequence r^l , along which $(\bar{X}^r, \bar{Y}^r, \bar{W}^r)$ converges uniformly to $(x_{m_l}, y_{m_l}, w_{m_l})$, which satisfies properties (9), at any $t \in [0, T_l]$. The limit trajectory $(x_{m_l}, y_{m_l}, w_{m_l})$ does not hit m_l in $[0, T_l]$. Subsequence $r^l = \{r_1^l, r_2^l, \dots\}$ is such that, w.p.1, for all sufficiently large r along the subsequence r^l , $(\bar{X}^r(t), \bar{Y}^r(t), \bar{W}^r(t)) = (\bar{X}_{m_l}^r(t), \bar{Y}_{m_l}^r(t), \bar{W}_{m_l}^r(t))$ at any $t \in [0, T_l]$. We consider a sequence $T_1, T_2, \dots, \rightarrow \infty$. We construct a subsequence r^* by using Cantor's diagonal process [16] from subsequences r^1, r^2, \dots ($r^1 \supseteq r^2 \supseteq \dots$) corresponding to T_1, T_2, \dots , respectively (i.e. $r_1^* = r_1^1, r_2^* = r_2^2, \dots$). Clearly, for this subsequence r^* , w.p.1, $(\bar{X}^r, \bar{Y}^r, \bar{W}^r)$ converges u.o.c. to (x, y, w) , which satisfies properties (9), at any regular point $t \in [0, \infty)$. \square

6 Proof of Theorem 2

We use the machinery of Common Quadratic Lyapunov Functions (CQLF) to approach the stability of fluid limits, i.e. their convergence to the unique equilibrium point $(0, 0, 0)$ [10][17].

System (10) is a switched linear system with $m = 2$. Namely, for $y \geq 0$,

$$\begin{cases} x'(t) = (-\gamma\beta)x + (\frac{1}{2}\gamma\alpha\mu - \epsilon)y + (-\frac{1}{2}\gamma\alpha\mu)w \\ y'(t) = (\beta)x + (-\frac{1}{2}\alpha\mu)y + (\frac{1}{2}\alpha\mu)w \\ w'(t) = (\beta)x + (-\frac{1}{2}(\alpha - 2)\mu)y + (\frac{1}{2}(\alpha - 2)\mu)w \end{cases} \quad (40)$$

and for $y < 0$,

$$\begin{cases} x'(t) = (-\gamma\beta)x + (-\frac{1}{2}\gamma\alpha\mu - \epsilon)y + (-\frac{1}{2}\gamma\alpha\mu)w \\ y'(t) = (\beta)x + (\frac{1}{2}\alpha\mu)y + (\frac{1}{2}\alpha\mu)w \\ w'(t) = (\beta)x + (\frac{1}{2}(\alpha - 2)\mu)y + (\frac{1}{2}(\alpha - 2)\mu)w. \end{cases} \quad (41)$$

We can rewrite the systems above as two linear time-invariant systems $u'(t) = A^+u(t)$ and $u'(t) = A^-u(t)$, where $u(t) = (x(t), y(t), w(t))^T$ and

$$A^+ = \begin{pmatrix} -\gamma\beta & \frac{1}{2}\gamma\alpha\mu - \epsilon & -\frac{1}{2}\gamma\alpha\mu \\ \beta & -\frac{1}{2}\alpha\mu & \frac{1}{2}\alpha\mu \\ \beta & -\frac{1}{2}(\alpha - 2)\mu & \frac{1}{2}(\alpha - 2)\mu \end{pmatrix} \quad (42)$$

and

$$A^- = \begin{pmatrix} -\gamma\beta & -\frac{1}{2}\gamma\alpha\mu - \epsilon & -\frac{1}{2}\gamma\alpha\mu \\ \beta & \frac{1}{2}\alpha\mu & \frac{1}{2}\alpha\mu \\ \beta & \frac{1}{2}(\alpha - 2)\mu & \frac{1}{2}(\alpha - 2)\mu \end{pmatrix}. \quad (43)$$

Lemma 2. Matrix A^+ in (42) is Hurwitz for all positive $\beta, \gamma, \mu, \epsilon$ and $\alpha \in (0, 1)$.

Proof. The characteristic equation of A^+ is $\det(A^+ - \lambda I) = 0$, which is equivalent to

$$\lambda^3 + (\beta\gamma + \mu)\lambda^2 + (\beta\epsilon + \beta\gamma\mu)\lambda + \beta\epsilon\mu = 0. \quad (44)$$

By Proposition 3, it suffices to verify that

$$\beta\gamma + \mu > 0, \quad \beta\epsilon + \beta\gamma\mu > 0, \quad \beta\epsilon\mu > 0, \quad (45)$$

and

$$(\beta\gamma + \mu)(\beta\epsilon + \beta\gamma\mu) - \beta\epsilon\mu = \beta^2\gamma^2\mu + \beta^2\gamma\epsilon + \beta\gamma\mu^2 > 0. \quad (46)$$

Conditions (45) and (46) are obviously true. \square

Lemma 3. Matrix A^- in (43) is Hurwitz for positive $\beta, \gamma, \mu, \epsilon$, and $\alpha \in (0, 1)$, satisfying

$$\left(\frac{\beta\gamma}{\mu} + (1 - \alpha)\right) \left(\frac{\gamma\mu}{\epsilon} + 1\right) > 1. \quad (47)$$

Proof. The characteristic equation of A^- is $\det(A^- - \lambda I) = 0$, which is equivalent to

$$\lambda^3 + (\beta\gamma + \mu(1 - \alpha))\lambda^2 + (\beta\epsilon + \beta\gamma\mu)\lambda + \beta\epsilon\mu = 0. \quad (48)$$

By Proposition 3, it suffices to verify that

$$\beta\gamma + \mu(1 - \alpha) > 0, \quad \beta\epsilon + \beta\gamma\mu > 0, \quad \beta\epsilon\mu > 0, \quad (49)$$

and

$$(\beta\gamma + \mu(1 - \alpha))(\beta\epsilon + \beta\gamma\mu) - \beta\epsilon\mu > 0 \text{ which is equivalent to } \left(\frac{\beta\gamma}{\mu} + (1 - \alpha)\right) \left(\frac{\gamma\mu}{\epsilon} + 1\right) > 1.$$

Conditions (49) are obviously true. \square

Lemma 4. For $\beta > 0, \mu > 0$ and $\alpha \in (0, 1)$, there exists a pair of $\gamma > 0$ and $\epsilon > 0$ satisfying conditions

$$\begin{cases} \frac{\beta\gamma^2}{4} < \epsilon < \frac{\beta\gamma^2}{2} \\ \epsilon > \frac{\beta\gamma^2}{2} - \left(\frac{\alpha\gamma\mu}{2} - \frac{(1-\alpha)\mu^2}{2\beta}\right) \\ \gamma > \frac{(1-\alpha)\mu}{\alpha\beta}. \end{cases} \quad (50)$$

Moreover, conditions (50) imply matrix A^- being Hurwitz.

Proof. For $\beta > 0, \mu > 0$ and $\alpha \in (0, 1)$, we have $\frac{(1-\alpha)\mu}{\alpha\beta} > 0$. Hence, we can always find a value of $\gamma > 0$ satisfying the third condition of (50). And from the third condition of (50), we have

$$\frac{\alpha\gamma\mu}{2} - \frac{(1-\alpha)\mu^2}{2\beta} > 0. \quad (51)$$

Hence, we can always find a value of $\epsilon > 0$ satisfying conditions

$$\begin{cases} \frac{\beta\gamma^2}{4} < \epsilon < \frac{\beta\gamma^2}{2} \\ \epsilon > \frac{\beta\gamma^2}{2} - \left(\frac{\alpha\gamma\mu}{2} - \frac{(1-\alpha)\mu^2}{2\beta}\right). \end{cases} \quad (52)$$

As shown in the proof of Lemma 3, matrix A^- when

$$(\beta\gamma + \mu(1 - \alpha))(\beta\epsilon + \beta\gamma\mu) - \beta\epsilon\mu = \beta^2\gamma\epsilon + \beta\epsilon\mu(1 - \alpha) + \beta^2\gamma^2\mu + \beta\gamma\mu^2(1 - \alpha) - \beta\epsilon\mu > 0.$$

For positive $\beta, \gamma, \mu, \epsilon$, and $\alpha \in (0, 1)$, the condition $\beta^2\gamma^2\mu - \beta\epsilon\mu > 0$ or, equivalently, $\epsilon < \gamma^2\beta$ implies condition (47). It means that, if $\epsilon < \gamma^2\beta$, then A^- is Hurwitz. But, the conditions (50) imply $\epsilon < \gamma^2\beta$. \square

Conclusion of the proof of Theorem 2. The characteristic equation of A^+A^- is

$$\lambda^3 - (\mu^2 - \alpha\mu^2 + \beta^2\gamma^2 - 2\beta\epsilon - \alpha\beta\gamma\mu)\lambda^2 + (\beta^2\epsilon^2 + \beta^2\gamma^2\mu^2 - 2\beta\epsilon\mu^2 + \alpha\beta\epsilon\mu^2)\lambda - \beta^2\epsilon^2\mu^2 = 0. \quad (53)$$

(Expression (53) is obtained with the help of MATLAB symbolic calculation.) By Proposition 4, if $\Delta < 0$, then the equation has one real root and two nonreal complex conjugate roots. It is well known that the determinant of a square matrix A^+A^- is the product of its eigenvalues. We have $\det(A^+A^-) = \lambda_1\lambda_2\lambda_3 = \beta^2\epsilon^2\mu^2 > 0$. Therefore, one of the roots must be a real positive. We see that it will suffice to show that $\Delta < 0$ to demonstrate A^+A^- could have no negative real eigenvalues. From (53), we have

$$\begin{cases} a = 1 \\ b = -(\mu^2 - \alpha\mu^2 + \beta^2\gamma^2 - 2\beta\epsilon - \alpha\beta\gamma\mu) \\ c = \beta^2\epsilon^2 + \beta^2\gamma^2\mu^2 - 2\beta\epsilon\mu^2 + \alpha\beta\epsilon\mu^2 \\ d = -\beta^2\epsilon^2\mu^2. \end{cases} \quad (54)$$

These a, b, c , and d are the coefficients of general cubic equation (4). From (5), we have

$$\Delta = 18bcd - 4b^3d + b^2c^2 - 4c^3 - 27d^2 = d((18c - 4b^2)b - 27d) + c^2(b^2 - 4c). \quad (55)$$

From (11), we have $c = \beta^2\epsilon^2 + \beta\mu^2(\beta\gamma^2 - 2\epsilon) + \alpha\beta\epsilon\mu^2 > 0$ (note that: $\beta\gamma^2 - 2\epsilon > 0$) and $d < 0$. Hence, to show that $\Delta < 0$ in equation (55), it will suffice to show

$$\begin{cases} b > 0 \\ b^2 - 4c < 0. \end{cases} \quad (56)$$

We will show that conditions (11) imply (56). We have

$$b = (\alpha - 1)\mu^2 - \beta^2\gamma^2 + 2\beta\epsilon + \alpha\beta\gamma\mu > (\alpha - 1)\mu^2 - \beta^2\gamma^2 + \alpha\beta\gamma\mu + \beta^2\gamma^2 - \alpha\beta\gamma\mu + (1 - \alpha)\mu^2 = 0$$

$$\left[\text{Note that } \epsilon > \frac{\beta\gamma^2}{2} - \left(\frac{\alpha\gamma\mu}{2} - \frac{(1 - \alpha)\mu^2}{2\beta} \right) \right],$$

and

$$\begin{aligned} b^2 - 4c &= \alpha^2\beta^2\gamma^2\mu^2 + 2\alpha^2\beta\gamma\mu^3 + \alpha^2\mu^4 - 2\alpha\beta^3\gamma^3\mu - 2\alpha\beta^2\gamma^2\mu^2 + \\ &+ 4\epsilon\alpha\beta^2\gamma\mu - 2\alpha\beta\gamma\mu^3 - 2\alpha\mu^4 + \beta^4\gamma^4 - 4\epsilon\beta^3\gamma^2 - 2\beta^2\gamma^2\mu^2 + 4\epsilon\beta\mu^2 + \mu^4 = \\ &= (\alpha - 1)^2\mu^4 + \beta\mu^2(\alpha^2\beta\gamma^2 - 2\alpha\beta\gamma^2 - 2\beta\gamma^2 + 4\epsilon) + 2\alpha\beta\gamma\mu^3(\alpha - 1) + \\ &\quad + \alpha\beta^2\gamma\mu(-2\beta\gamma^2 + 4\epsilon) + \beta^3\gamma^2(\beta\gamma^2 - 4\epsilon) \stackrel{(a)}{<} \\ &< (\alpha - 1)^2\mu^4 + \beta\mu^2(\alpha^2\beta\gamma^2 - 2\alpha\beta\gamma^2 - 2\beta\gamma^2 + 4\epsilon) + 2\alpha\beta\gamma\mu^3(\alpha - 1) = \\ &= (\alpha - 1)\mu^3((\alpha - 1)\mu + \alpha\beta\gamma) + \alpha\beta\gamma\mu^3(\alpha - 1) + \beta\mu^2(\alpha\beta\gamma^2(\alpha - 2) - 2(\beta\gamma^2 - 2\epsilon)) \stackrel{(b)}{<} \\ &< (\alpha - 1)\mu^3((\alpha - 1)\mu + \alpha\beta\gamma). \end{aligned}$$

(where in (a) and (b) we use the facts that $\frac{\beta\gamma^2}{4} < \epsilon < \frac{\beta\gamma^2}{2} \Leftrightarrow \beta\gamma^2 - 4\epsilon < 0 < \beta\gamma^2 - 2\epsilon$).

From conditions (11), we have $(\alpha - 1)\mu + \alpha\beta\gamma > (\alpha - 1)\mu + (1 - \alpha)\mu = 0$. Note that $\gamma > \frac{(1-\alpha)\mu}{\alpha\beta} \Leftrightarrow \alpha\beta\gamma > (1 - \alpha)\mu$. Therefore, we have $b > 0$ and $b^2 - 4c < 0$. Hence, A^+A^- has no negative real eigenvalues under conditions (11).

By Lemma 2, A^+ is Hurwitz for all positive $\beta, \gamma, \mu, \epsilon$ and $\alpha \in (0, 1)$. By Lemma 4, A^- is Hurwitz under conditions (11). It is easy to verify that the difference $A^+ - A^-$ has rank one. A^+A^- has no negative real eigenvalues under conditions (11). Hence, by Proposition 2, $u'(t) = A^+u(t)$ and $u'(t) = A^-u(t)$ have a CQLF. Therefore, by Proposition 1, the system (10) is exponentially stable under conditions (11). This completes the proof. \square

As a useful corollary of Lemma 3, we have the following fact.

Corollary 1. *If matrix A^- in (43) is Hurwitz for some positive $\beta, \gamma, \mu, \epsilon$, and $\alpha \in (0, 1)$, then it remains Hurwitz if α is replaced by any $0 < \alpha_0 \leq \alpha$.*

Proof. From condition (47), for any $\alpha_0 \in (0, \alpha]$, we have

$$\left(\frac{\beta\gamma}{\mu} + (1 - \alpha_0)\right) \left(\frac{\gamma\mu}{\epsilon} + 1\right) \geq \left(\frac{\beta\gamma}{\mu} + (1 - \alpha)\right) \left(\frac{\gamma\mu}{\epsilon} + 1\right) > 1. \quad (57)$$

Application of Lemma 3 completes the proof. \square

7 Numerical examples

In this section, we present some numerical examples to show the good performance of the scheme. Later, we also provide some conjectures based on a variety of simulations.

Example 1. We use the following set of parameters satisfying the conditions (11):

$$\Lambda = 1000, \alpha = 0.7, \beta = 1, \mu = 1, \gamma = 2, \epsilon = 1.5$$

We consider two initial conditions: (i) $(X(0), Y(0), Z(0)) = (0, 0, 0)$; (ii) $(X(0), Y(0), Z(0)) = (0, -1000, 0)$ (Figure 2). The red line of the figure is the fluid limit and the blue line of the figure is the simulation experiment. The results suggest the global stability of the system.

Example 2. We use another set of parameters satisfying the conditions (11):

$$\Lambda = 1000, \alpha = 0.5, \beta = 3, \mu = 2, \gamma = 1, \epsilon = 1.4$$

We consider two initial conditions: (i) $(X(0), Y(0), Z(0)) = (2000, 0, 1000)$; (ii) $(X(0), Y(0), Z(0)) = (0, 2000, 0)$ (Figure 3). The results suggest the global stability of the system even if the boundary on x is hit (on the left of the Figure 3).

Example 3. We use 4 another sets of parameters (with different values of α), which do not satisfy the conditions (11), but satisfy A^- Hurwitz condition (47):

$$\Lambda = 1000, \beta = 1, \mu = 2, \gamma = 2, \epsilon = 0.2$$

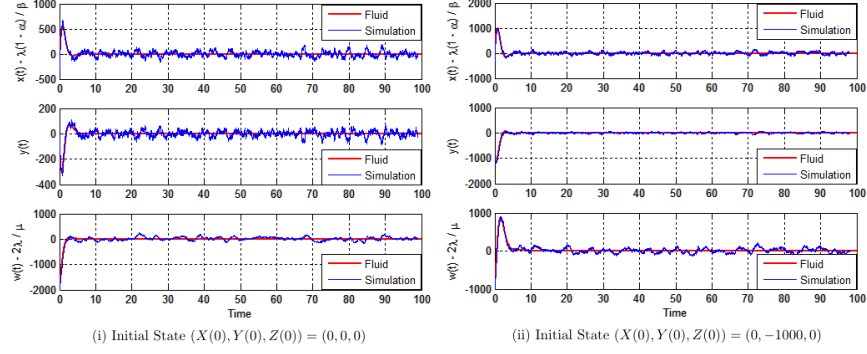


Figure 2: Comparison of fluid approximations with simulations in Example 1

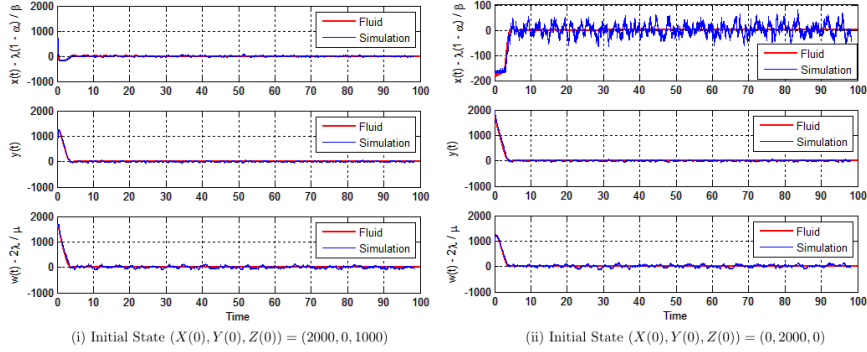


Figure 3: Comparison of fluid approximations with simulations in Example 2

We consider an initial condition $(X(0), Y(0), Z(0)) = (0, 1000, 500)$ with 4 different values of α ($\alpha_1 = 0.1$, $\alpha_2 = 0.4$, $\alpha_3 = 0.6$, and $\alpha_4 = 0.9$) (Figure 4). The results suggest the global stability of the system.

Example 4. We use another sets of parameters, which do not satisfy A^- Hurwitz condition (47):

- (i) $\Lambda = 1000$, $\alpha = 0.5$, $\beta = 0.05$, $\mu = 0.5$, $\gamma = 1$, $\epsilon = 1$, and
- (ii) $\Lambda = 1000$, $\alpha = 0.9$, $\beta = 0.05$, $\mu = 0.5$, $\gamma = 1$, $\epsilon = 1$

We consider an initial condition $(X(0), Y(0), Z(0)) = (500, 1000, 500)$ (Figure 5). On the left of the figure shows that the system is globally stable, but unstable on the right of the figure.

The results of Example 1 and Example 2 suggest that the system is globally stable under the conditions (11) even if the boundary on x is hit or not. The results of Example 3 suggest that the system is locally and globally stable when A^- is Hurwitz even if the conditions (11) are not satisfied. The results of Example 4 show that the system can be globally stable under some set of parameters, but unstable under some different set of parameters when A^- is not Hurwitz. These numerical results suggest the following conjectures:

Conjecture 1. The system is globally stable if it is locally stable.

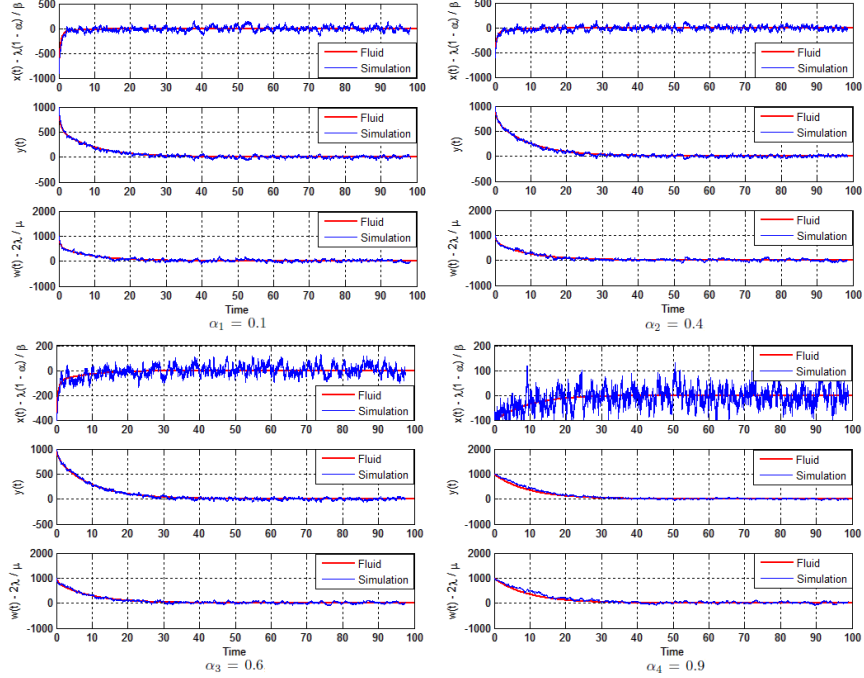


Figure 4: Comparison of fluid approximations with simulations in Example 3

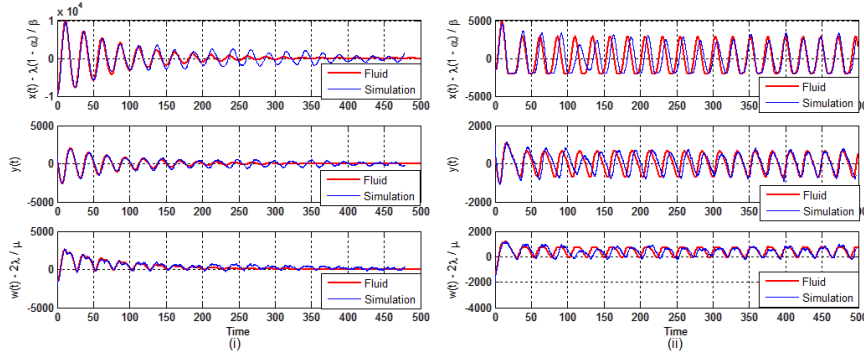


Figure 5: Comparison of fluid approximations with simulations in Example 4

Conjecture 2. If A^- is Hurwitz, it is sufficient for local stability of the system. (A^+ is always Hurwitz in our case.)

8 Conclusions

In this paper, we study a feedback-based agent invitation scheme for a model with randomly behaving agents. This model is motivated by a variety of existing and emerging applications. The focus of the paper is on the stability properties of the system fluid limits, arising as asymptotic limits of the system process, when the system scale (customer arrival rate) grows to infinity. The dynamic system, describing the behavior of fluid limit trajectories has a very complex structure – it is a switched linear system, which in addition has a reflecting boundary. We derived some sufficient

local stability conditions, using the machinery of switched linear systems and Common Quadratic Lyapunov Functions. Our simulation and numerical experiments show good overall performance of the feedback scheme, when the local stability conditions hold. They also suggest that the local stability is in fact sufficient for the global stability of fluid limits. Verifying these conjectures, as well as expanding the sufficient local stability conditions, is an interesting subject for future research. Further generalizations of the agent invitation model are also of interest from both theoretical and practical points of view.

References

- [1] Z. Aksin, M. Armony, and V. Mehrotra. The modern call center: a multi-disciplinary perspective on operations management research. *Production and Operations Management*, 16(6):665–688, 2007.
- [2] American Telemedicine Association. *Core Operational Guidelines for Telehealth Services Involving Provider-Patient Interactions*, 2014. <http://www.americantelemed.org/docs/default-source/standards/core-operational-guidelines-for-telehealth-services.pdf?sfvrsn=6>.
- [3] S. Bengtson. Generating better results with crowdsourcing: Leverage a network of high-quality professionals for customer service. *White paper*, 2014.
- [4] S. N. Ethier and T. G. Kurtz. *Markov Processes: Characterization and Convergence*. Wiley, 1986.
- [5] P. Formisano. Flexibility for changing business needs: Improve customer service and drive more revenue with a virtual crowdsourcing solution. *White paper*, 2014.
- [6] I. Gurvich and A. Ward. On the dynamic control of matching queues. *Stochastic Systems*, 4(2):479–523, 2014.
- [7] J. P. Hespanha. Uniform stability of switched linear systems: extensions of LaSalle’s invariance principle. *IEEE Transactions on Automatic Control*, 49(4):470–482, 2004.
- [8] R. S. Irving. *Integers, Polynomials, and Rings*. Undergraduate Texts in Mathematics. Springer, New York, USA, 2004 edition, 2004.
- [9] B. R. K. Kashyap. The double-ended queue with bulk service and limited waiting space. *Operations Research*, 14(5):822–834, 1966.
- [10] H. Lin and P. J. Antsaklis. Stability and stabilizability of switched linear systems: A survey of recent results. *IEEE Transactions on Automatic Control*, 54(2):308–322, 2009.
- [11] X. Liu, Q. Gong, and V. G. Kulkarni. Diffusion models for doubly-ended queues with renewal arrival processes. Forthcoming in *Stochastic Systems*. DOI: 10.1214/13-SSY113, 2014.
- [12] S. McGee-Smith. Why companies are choosing to deploy the LiveOps cloud-based contact center, 2010. http://www.liveops.com/sites/default/files/uploads/lo_wp_mcgee-smith_analytics.pdf.
- [13] G. Pang and A. Stolyar. A service system with on-demand agent invitations. *Queueing Systems*, 2015. <http://arxiv.org/pdf/1409.7380v2.pdf>.

- [14] G. Pang, R. Talreja, and W. Whitt. Martingale proofs of many-server heavy-traffic limits for Markovian queues. *Probability Surveys*, 4:193–267, 2007.
- [15] L. Pontryagin. *Ordinary Differential Equations*. Adiwes international series in mathematics. Addison-Wesley, 1962.
- [16] W. Rudin. *Principles of Mathematical Analysis*. McGraw-Hill, New York, 1976.
- [17] R. Shorten, F. Wirth, O. Mason, K. Wulff, and C. King. Stability criteria for switched and hybrid systems. *SIAM Review*, 49(4):545–592, 2007.
- [18] A. Stolyar, M. Reiman, N. Korolev, V. Mezhibovsky, and H. Ristock. Pacing in knowledge worker engagement, 2010. *United States Patent Application 20100266116-A1*.